

It's important to evaluate the performance of the learned hypothesis (in many cases, hypothesis evaluation is even an integral component of the learning method itself). When evaluating a learned hypothesis we are interested in estimating the accuracy with which it will classify future instances.

Accuracy can be expressed in terms of error:

$$\text{accuracy}(h) = 1 - \text{error}(h)$$

- true error : error rate of the hypothesis over the entire unknown distribution Δ of examples

The true error of hypothesis h with respect to target function f and distribution Δ is the probability that h will misclassify an instance drawn at random according to Δ

$$\text{distribution} \quad \text{error}_\Delta(h) = \Pr_{x \in \Delta} [f(x) \neq h(x)]$$

- sample error : error rate of the hypothesis over the sample of data that is available

The sample error of h with respect to target function f and data sample S is the proportion of examples that h misclassifies

$$\text{error}_S(h) = \frac{1}{n} \sum_{x \in S} \delta(f(x) \neq h(x))$$

where $\delta(f(x) \neq h(x))$ is 1 if $f(x) \neq h(x)$ and 0 otherwise

we can't compute the true error (we don't know the distribution). We can compute the sample error, but only on a small data sample

How well does the sample error ($\text{error}_S(h)$) estimate the true error ($\text{error}_\Delta(h)$)?

The learning system should be as accurate as possible $\rightarrow \forall x \notin S$

Overfitting $\left\{ \begin{array}{l} \text{accuracy}_S(h) \text{ is very high} \\ \text{accuracy}_S(h) \text{ is very low} \end{array} \right\}$ hypothesis $h \in H$ overfits training data if there is an alternative hypothesis $h' \in H$ such that $\text{error}_S(h) < \text{error}_S(h')$ and $\text{error}_S(h) > \text{error}_S(h')$

How to compute the sample error

- partition the dataset Δ ($\Delta = T \cup S$, $T \cap S = \emptyset$, $|T| = 2/3 |\Delta|$); the split must be done in a random way
- compute a hypothesis h using training set T
- evaluate $\text{error}_S(h) = \frac{1}{n} \sum_{x \in S} \delta(f(x) \neq h(x))$

The sample error is an unbiased estimator for $\text{error}_\Delta(h)$

Trade off between training and testing

- having more samples for training and less for testing improves the performance of the model but $\text{error}_S(h)$ does not approximate well $\text{error}_\Delta(h)$
- having more samples for evaluation and less for training reduces variance of estimation: $\text{error}_S(h)$ approximates well $\text{error}_\Delta(h)$ but this value may be not satisfactory.
In general, for medium sized datasets: 2/3 is used for training and 1/3 for testing

True comparison between two hypothesis: $d = \text{error}_\Delta(h_1) - \text{error}_\Delta(h_2)$

Estimator of the comparison between two hypothesis: $\hat{d} = \text{error}_{S_1}(h_1) - \text{error}_{S_2}(h_2)$

$\hookrightarrow \hat{d}$ is an unbiased estimator for d iff h_1, h_2, S_1, S_2 are independent from each other

h is the solution of learning algorithm L when using a training set $T \longrightarrow h = L(T)$

Often we are interested in comparing the performance of two learning algorithms L_a and L_b , rather than two specific hypotheses.

$\text{error}_{S_i}(h)$ is the result of only one experiment and it is not informative on the learning algorithm. We can perform many experiments and compute $\text{error}_{S_i}(h)$ for different independent sample data S_i : this is called **K-fold Cross Validation** method.

K-fold Cross Validation

- partition data set D into K disjoint sets S_1, S_2, \dots, S_K ($|S_i| > 30$)
- for $i = 1, \dots, K$ do
 - use S_i as test set, and the remaining data as training set T_i
 - $T_i \leftarrow \{D - S_i\}$
 - $h_i \leftarrow L(T_i)$
 - $\delta_i \leftarrow \text{error}_{S_i}(h_i)$
- return $\text{error}_{L,D} = \frac{1}{K} \sum_{i=1}^K \delta_i$

error is the mean value
of all errors computed

$$\text{accuracy}_{L,D} = 1 - \text{error}_{L,D}$$

Repeating the process several times we will obtain different hypotheses and different values of $\text{error}_{S_i}(h_i)$

Comparing learning algorithms L_a and L_b using K-fold Cross Validation

- partition data set D into K disjoint sets S_1, S_2, \dots, S_K ($|S_i| > 30$)
- for $i = 1, \dots, K$ do
 - use S_i as test set, and the remaining data as training set T_i
 - $T_i \leftarrow \{D - S_i\}$
 - $h_{A,i} \leftarrow L_a(T_i)$
 - $h_{B,i} \leftarrow L_b(T_i)$
 - $\delta_i \leftarrow \text{error}_{S_i}(h_{A,i}) - \text{error}_{S_i}(h_{B,i})$
- return $\bar{\delta} = \frac{1}{K} \sum_{i=1}^K \delta_i$

error is the mean value
of all errors computed

$$\text{accuracy}_{L,D} = 1 - \text{error}_{L,D}$$

if $\bar{\delta} < 0$ then L_a is better than L_b ($\sim \text{error}_{S_i}(h_{A,i}) - \text{error}_{S_i}(h_{B,i}) < 0 \longrightarrow \text{error}_{S_i}(h_{A,i}) < \text{error}_{S_i}(h_{B,i})$)
 if $\bar{\delta} > 0$ then L_b is better than L_a ($\sim \text{error}_{S_i}(h_{B,i}) - \text{error}_{S_i}(h_{A,i}) < 0 \longrightarrow \text{error}_{S_i}(h_{B,i}) < \text{error}_{S_i}(h_{A,i})$)

Accuracy only is not enough to assess the performance of a classification method.

Unbalanced data sets are very common in problems related to anomaly detection (e.g. malware analysis, fraud detection, medical tests, ...)

Performance metrics

→ Concept learning = classification problem with two classes

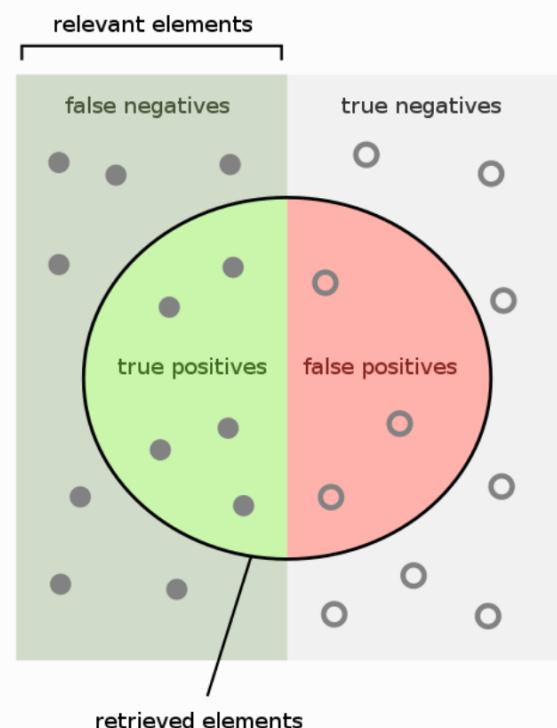
		Predicted class	
True class	YES	NO	
YES	True positive	False negative	
NO	False positive	True negative	

Precision = ability to avoid false positives = $\frac{|TP|}{|TP+FP|}$
 = ability to correctly classify positive samples (avoid FP)

Recall = ability to avoid false negatives = $\frac{|TP|}{|TP+FN|}$
 precision can be seen as a measure of fidelity (avoid FN)

Precision is the number of true positives compared to the total number of positives identified (true positives + false positives)

Recall is the number of true positives compared to all the existing positives (true positives + false negatives)



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{green}}{\text{green} + \text{red}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{green}}{\text{green} + \text{black}}$$

→ Classification problem with many (more than two) classes

In a classification problem with many classes, we can compute how many times an instance of class C_i is classified in class C_j .

	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8
C_1	■							
C_2		■						
C_3			■					
C_4				■				
C_5					■			
C_6						■		
C_7							■	
C_8								■

confusion matrix

Other performance measures

$$F_1\text{-score} = 2(\text{Precision} \cdot \text{Recall}) / (\text{Precision} + \text{Recall})$$

$$\text{True Positive Rate TPR (Sensitivity)} = \frac{TP}{P} = \frac{TP}{TP+FN}$$

$$\text{True Negative Rate TNR (Specificity)} = \frac{TN}{N} = \frac{TN}{TN+FP}$$

$$\text{False Positive Rate FPR} = \frac{FP}{N} = \frac{FP}{TN+FP}$$

$$\text{False Negative Rate FNR} = \frac{FN}{P} = \frac{FN}{TP+FN}$$

Main diagonal contains the accuracy for each class.
 Outside the diagonal, the errors.
 Cell C_{ij} contains the times class C_i is classified as class C_j .

e.g. $C_{1,2}$ gives us information on how many times an instance of class C_1 is classified as class C_2